



ARL-TR-8304 • FEB 2018



Machine-Learning Techniques for the Determination of Attrition of Forces Due to Atmospheric Conditions

by Yasmina R Raby

Approved for public release; distribution is unlimited.

NOTICES

Disclaimers

The findings in this report are not to be construed as an official Department of the Army position unless so designated by other authorized documents.

Citation of manufacturer's or trade names does not constitute an official endorsement or approval of the use thereof.

Destroy this report when it is no longer needed. Do not return it to the originator.



Machine-Learning Techniques for the Determination of Attrition of Forces Due to Atmospheric Conditions

by Yasmina R Raby

Computational and Information Sciences Directorate, ARL

REPORT DOCUMENTATION PAGE				Form Approved OMB No. 0704-0188	
<p>Public reporting burden for this collection of information is estimated to average 1 hour per response, including the time for reviewing instructions, searching existing data sources, gathering and maintaining the data needed, and completing and reviewing the collection information. Send comments regarding this burden estimate or any other aspect of this collection of information, including suggestions for reducing the burden, to Department of Defense, Washington Headquarters Services, Directorate for Information Operations and Reports (0704-0188), 1215 Jefferson Davis Highway, Suite 1204, Arlington, VA 22202-4302. Respondents should be aware that notwithstanding any other provision of law, no person shall be subject to any penalty for failing to comply with a collection of information if it does not display a currently valid OMB control number.</p> <p>PLEASE DO NOT RETURN YOUR FORM TO THE ABOVE ADDRESS.</p>					
1. REPORT DATE (DD-MM-YYYY) February 2018		2. REPORT TYPE Technical Report		3. DATES COVERED (From - To) October 2015–September 2017	
4. TITLE AND SUBTITLE Machine-Learning Techniques for the Determination of Attrition of Forces Due to Atmospheric Conditions				5a. CONTRACT NUMBER	
				5b. GRANT NUMBER	
				5c. PROGRAM ELEMENT NUMBER	
6. AUTHOR(S) Yasmina R Raby				5d. PROJECT NUMBER	
				5e. TASK NUMBER	
				5f. WORK UNIT NUMBER	
7. PERFORMING ORGANIZATION NAME(S) AND ADDRESS(ES) US Army Research Laboratory Computational and Information Sciences Directorate (ATTN: RDRL-CIE-D) White Sands Missile Range, NM 88002-5513				8. PERFORMING ORGANIZATION REPORT NUMBER ARL-TR-8304	
9. SPONSORING/MONITORING AGENCY NAME(S) AND ADDRESS(ES)				10. SPONSOR/MONITOR'S ACRONYM(S)	
				11. SPONSOR/MONITOR'S REPORT NUMBER(S)	
12. DISTRIBUTION/AVAILABILITY STATEMENT Approved for public release; distribution is unlimited.					
13. SUPPLEMENTARY NOTES					
14. ABSTRACT <p>This report documents the findings of an attempt to model the attrition of forces due to atmospheric conditions. Machine-learning techniques, primarily the random forest algorithm, were used to explore the possibility of a correlation between aircraft incidents in the National Transportation Safety Board database and meteorological conditions. If a strong correlation could be found, it could be used to derive a model to predict aircraft incidents and become part of a decision support tool for mission planning purposes. While the random forest algorithm was able to discover some consistent predictors across a variety of data sets while classifying aircraft incidents related to weather, there were some concerns regarding the error rate in the final result of the classification process. This report documents the efforts to define a model and provide lessons learned toward future attempts to refine the results and generate a model that addresses the attrition of forces due to atmospheric conditions using machine-learning techniques.</p>					
15. SUBJECT TERMS machine learning, decision aids, aircraft, attrition, random forest					
16. SECURITY CLASSIFICATION OF:			17. LIMITATION OF ABSTRACT UU	18. NUMBER OF PAGES 38	19a. NAME OF RESPONSIBLE PERSON Yasmina Raby
a. REPORT Unclassified	b. ABSTRACT Unclassified	c. THIS PAGE Unclassified			19b. TELEPHONE NUMBER (Include area code) (575) 678-1761

Contents

List of Figures	iv
List of Tables	v
Acknowledgments	vi
Summary	vii
1. Introduction	1
2. Methods, Assumptions, and Procedures	2
2.1 Data Collection	2
2.2 Random Forest Implementation	3
2.3 Random Forest Output	4
3. Results and Discussion	7
4. Conclusions	10
5. References	11
Appendix. Database Queries and Random Forest Output	13
List of Symbols, Abbreviations, and Acronyms	27
Distribution List	28

List of Figures

Fig. 1	Spreadsheet containing aircraft data from NTSB Access database ..	2
Fig. 2	Data input for random forest	4
Fig. 3	Example variable importance plot.....	6
Fig. 4	OOB estimate of error rates for airplane data sets	8
Fig. 5	OOB estimate of error rates for helicopter data sets	10
Fig. A-1	Variable importance plot: airplanes, ntree 500, no gusts	15
Fig. A-2	Variable importance plot: airplanes, ntree 400, no gusts	16
Fig. A-3	Variable importance plot: airplanes, ntree 500, gusts	17
Fig. A-4	Variable importance plot:-airplanes, ntree 400, gusts.....	19
Fig. A-5	Variable importance plot: airplanes, ntree 500, gusts, no altimeter.....	20
Fig. A-6	Variable importance measures: helicopter, ntree 500, no gusts	21
Fig. A-7	Variable importance plot: helicopter, ntree 400, no gusts.....	23
Fig. A-8	Variable importance plot: helicopter, gusts, ntree 500.....	24
Fig. A-9	Variable importance plot: helicopter, gusts, ntree 400.....	25
Fig. A-10	Variable importance plot for helicopter, no gusts, no altimeter, ntree 500	26

List of Tables

Table 1	Example confusion matrix	5
Table 2	Example variable importance measures	6
Table 3	Summary of results on airplane data set.....	8
Table 4	Summary of results on helicopter data set.....	9
Table A-1	Confusion matrix airplanes, ntree 500, no gusts	14
Table A-2	Variable importance measures: airplanes, ntree 500, no gusts	15
Table A-3	Confusion matrix for airplanes, ntree 400, no gusts	16
Table A-4	Variable importance measures: airplanes, ntree 400, no gusts	16
Table A-5	Confusion matrix airplanes, ntree 500, gusts	17
Table A-6	Variable importance measures: airplanes, ntree 500, gusts	18
Table A-7	Confusion matrix, airplanes, ntree 400, gusts	18
Table A-8	Variable importance measures: airplanes, ntree 400, gusts	19
Table A-9	Confusion matrix, airplanes, ntree 500, no gusts, no altimeter	20
Table A-10	Variable importance measures: airplanes, ntree 500, no gusts, no altimeter.....	20
Table A-11	Confusion matrix, helicopters, ntree 500, no gusts	21
Table A-12	Variable importance plot: helicopter, ntree 500, gusts.....	22
Table A-13	Confusion matrix for helicopter, no gusts, ntree 400.....	22
Table A-14	Variable importance measures: helicopter, ntree 400, no gusts	22
Table A-15	Confusion matrix for helicopter, gusts, ntree 500	23
Table A-16	Importance measures for helicopter, gusts, ntree 500	24
Table A-17	Confusion matrix for helicopter, gusts, ntree 400	25
Table A-18	Importance measures for helicopter, gusts, ntree 400	25
Table A-19	Confusion matrix for helicopter, no gusts, no altimeter, ntree 500	26
Table A-20	Variable importance measures for helicopter, no gusts, no altimeter, ntree 500	26

Acknowledgments

The author thanks Mr Jeffrey O Johnson for his mentorship and providing the foundational concepts in his Human Vulnerability Forecasting project, and COL John R Olson (Ret) for the initial random forest algorithm and implementation.

Summary

This report documents the findings of an attempt to model the attrition of forces due to atmospheric conditions. Machine-learning techniques, primarily the random forest algorithm, were used to explore the possibility of a correlation between aircraft incidents in the National Transportation Safety Board database and meteorological conditions. If a strong correlation could be found, it could be used to derive a model to predict aircraft incidents and become part of a decision support tool for mission planning purposes.

While the random forest algorithm was able to discover some consistent predictors across a variety of data sets while classifying aircraft incidents related to weather, there were some concerns regarding the error rate in the final result of the classification process. This report documents the efforts to define a model and provide lessons learned toward future attempts to refine the results and generate a model that addresses the attrition of forces due to atmospheric conditions using machine-learning techniques.

INTENTIONALLY LEFT BLANK.

1. Introduction

Each year, the United States faces billions of dollars in damages and numerous deaths due to weather events.¹ Specific to Army operations, degraded visual environments have been the primary contributing factor to a majority of Army aviation mishaps during the past decade.² A known challenge to Department of Defense operations, quantifying potential losses due to combinations of geographic conditions and atmospheric effects in theater is critical to inform decision makers conducting combat and noncombat operations.

This issue became relevant in October 2015, when a Joint Land Attack Cruise Missile Defense Elevated Netted Sensor System (JLENS) aerostat broke loose from Aberdeen Proving Grounds and was widely reported by the media as it traveled freely, dragging 6700 ft of tether for 3 h, causing power outages as the tether damaged power lines. The blimp finally crashed in Pennsylvania, and the future of the program was put in jeopardy following the incident.³ This study attempts to resolve the issue of predicting what conditions would have resulted in the JLENS aerostat breaking loose and preventing such disasters in the future.

In general, this project intended to determine what atmospheric conditions contribute to aircraft losses while generating a model that will allow projected losses to be estimated given forecast data as input. Such information can be provided at the tactical operations center level to inform unit commanders, potentially improving mission success and unit survivability.

Initially, the focus for this project was to perform a feasibility study on the availability of data involved in events resulting in the loss of aerostats. However, insufficient data existed to create a significant determination for this study. Instead, a sizable database involving aircraft accidents from the National Transportation Safety Board (NTSB) was selected as a proof of concept due to its vast number of data points.

While this report does note some trends associated with temperature and dew point, there are some concerns about the error rates related to classification, and thus a model was not developed from the results. However, this technical report will serve as discussion of lessons learned during this process and as documentation of efforts to project aircraft losses using machine-learning techniques.

2. Methods, Assumptions, and Procedures

2.1 Data Collection

To obtain a sizable number of data points, the NTSB's Aviation and Accident Database was used as the data source for civilian aviation incidents. The NTSB provides this data in the form of Access databases, which can be broken down by date, type/make/model of aircraft involved, injuries, and so on.⁴ The database included incidents from January 1982 through October 2015. To narrow down the data sets, queries were formatted to create separate data sets for helicopters and airplanes, while selectively requesting the event IDs, descriptions of events, light conditions, temperature, dew point, wind direction, wind velocity, gusts, and altimeter. Queries developed for the collection of these data sets are presented in the Appendix.

To run the random-forest classification algorithm on this data, we had to distinguish a "response" variable that categorizes the event incident as occurring due to weather/environmental conditions or not. The "finding_description" field was searched for weather-related keywords using an Excel statement, and each record was marked as a response "x" for weather-related incidents and "N" for nonweather-related incidents.

An example of the resulting data is shown in Fig. 1:

	A	B	C	D	E	F	G	H	I
1	ev_id	light_conc	wx_temp	wx_dew	wind_dir	wind_vel	altimeter	Response	finding_description
2	20080122X00087	DAYL	19	16	70	4	30.31	N	Aircraft-Aircraft power plant-Engine (turbine/turboprop)-Turbine section-Fatigue/wear/corrosion - C
3	20080201X00130	NITE	13	8	100	3	30.2	N	Aircraft-Aircraft oper/perf/capability-Performance/control parameters-Altitude-Not attained/maintained - C
4	20080210X00162	DAYL	23	13	90	9	30.17	x	Environmental issues-Conditions/weather/phenomena-Temp/humidity/pressure-Conductive to carburetor icing-Effect on equipment - C
5	20080222X00227	NDRK	23	21	310	8	29.81	N	Aircraft-Aircraft oper/perf/capability-Performance/control parameters-(general)-Not attained/maintained - C

Fig. 1 Spreadsheet containing aircraft data from NTSB Access database

Data were broken up into several different categories to observe the impacts of each parameter on the results. For both airplane and helicopter records, the data sets were divided in the following ways:

- Gusts included or excluded
- 400 trees or 500 trees
- Gusts excluded; no altimeter and/or light conditions parameter; 500 trees

Gusts were excluded from the data sets due to gust data being rarely reported, and this results in an incomplete and thus unusable record. By excluding the gust data as a requirement, we kept an additional 4452 complete records for airplanes and 500 additional helicopter records. Additionally, we compared the results of using 500 trees versus 100 fewer trees. Based on some of the results of the previous entries, a final data set was created in an attempt to remove a parameter that could have added nonweather-related parameters in the data set. Thus, altimeter was removed while keeping as many records possible by excluding gusts.

For the helicopter data set, light conditions seemed to become highlighted as a most valuable predictor (MVP). While this is an important factor that should be considered in a model, we wanted to give other weather-related parameters an opportunity to become an MVP, and tried a final data set that included no gusts for the highest quantity of records, no altimeter, and no light condition.

2.2 Random Forest Implementation

This project uses the ‘randomForest’ R statistics-oriented language package, originally developed by Leo Breiman and Adele Cutler and ported for R by Andy Liaw and Matthew Wiener.⁵ The R source code that uses the randomForest package was written by COL John R Olson (Ret) and used in collaboration with the Human Vulnerability Forecasting project led by Jeffrey O Johnson.⁶ The randomForest implementation was chosen for its renowned accuracy and applicability for classification and regression problems that this topic relies on.⁷

For each of the runs, the randomForest implementation was modified to use different parameters and data sets. Each data set was modified to support each of the parameters drawn from the NTSB database and the desired data set required. Data were cropped only to include specific parameters and the response variable, as shown in Fig. 2.

light_cond	wx_temp	wx_dew_pt	wind_dir_deg	wind_vel_kts	altimeter	sky_ceil_ht	sky_nonceil_ht	Response
NDRK	20	13	30	23	30.13	5500	4100	x
DAYL	1	1	240	3	30.35	100	100	x
NITE	18	14	330	10	30	10000	3500	N
DAYL	15	4	190	14	29.92	5500	4800	x
DAYL	15	12	140	6	30.27	500	500	N
DAYL	17	1	260	9	29.99	9000	7000	x
NITE	1	-9	210	4	29.98	9500	9500	x
DAYL	28	18	70	15	30.08	7500	3300	x

Fig. 2 Data input for random forest

In each random forest execution, the importance variable was set to “TRUE”, allowing us to obtain a variable importance measure, as defined by the random forest package.⁵ The variable importance plot and importance values were then requested for each data set.

2.3 Random Forest Output

For each data set, our classification random-forest implementation output the following information:

- Type of random forest
- Number of trees
- Number of variables tried at each split
- Out-of-bag error (OOB) estimate
- Confusion matrix
- Variable importance plot
- Table of variable importance data

For each of the cases, we used classification models, and the number of variables tried remained at 2, meaning we used 2 random features to grow a single tree. We attempted to use the standard number of trees (400–500) to give us some variation of the impacts of using more or less trees on the number of parameters involved. A large number of trees was chosen to increase the likelihood of an accurate error estimate.⁸

The OOB error provides an estimate of errors from the random-forest algorithm and allows a glance of the performance of the algorithm on that particular data set. When the algorithm starts, the bootstrap training sample omits a third of the cases. These OOB cases are put back into their tree, which is tested to measure the

decrease in estimated margins and indicates whether or not the variable is important.⁹ This allows observers to decide whether to give credence to the MVPs that have been selected by the algorithm.

The confusion matrix provides another means for determining the veracity of the result by allowing us to see the predicted response against the actual response from the data set. The confusion matrices listed in the Appendix have columns listed for the predicted events (“N” for nonweather events and “x” for weather events) and the error rate for that class. The rows are labeled for the actual occurrence of those events. Thus, for every row–column combination we determined how many predictions were accurately classified versus how many were incorrectly classified. A confusion matrix describes the actual versus predicted values. The N row is the actual number of nonweather events, and the x row is the actual number of weather-related events. The N column is the predicted nonweather events, while the x column is the predicted weather events. Thus, the principal diagonal (N,N and x,x) represents the number of accurately predicted nonweather and weather events.

In Table 1, the N column demonstrates that of the 5589 records there were 5066 (3963 + 1103) total nonweather events, and the x column shows there were 523 (260 + 263) weather-related events. Of the 5066 nonweather related events, 3963 were predicted accurately. Similarly, of the 523 weather-related events, 263 were predicted accurately.

Table 1 Example confusion matrix

	N	x	class.error
N	3963	260	0.06156761
x	1103	263	0.80746706

The variable importance plot describes which parameters were determined to be the most successful among them. For each variable importance plot, the variable importance measure has also been provided, and more details about the plotted values are described in Fig. 3. The results of this plot should be weighed against the error estimates previously mentioned.

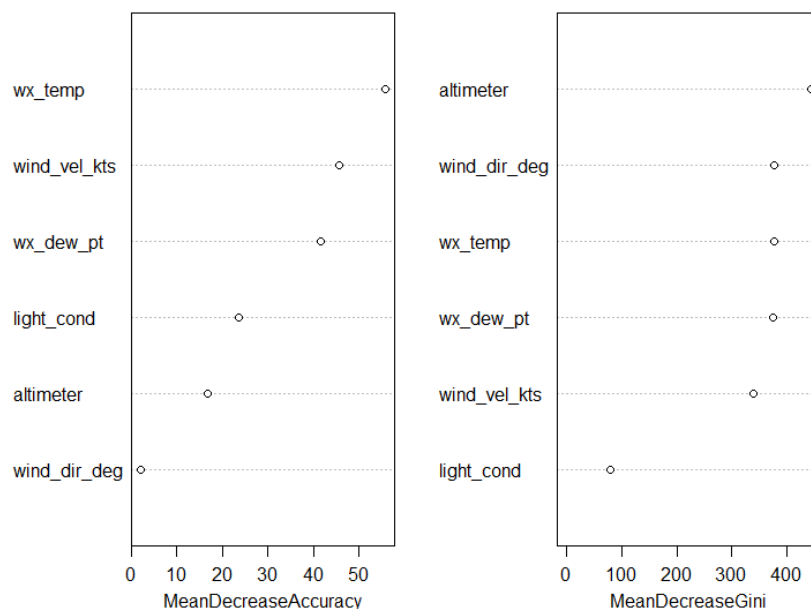


Fig. 3 Example variable importance plot

The variable importance measures in Table 2 provide detailed data on each of the predictors and the measure of importance as predictors. They do this by determining the increase of the OOB error estimate changed when a particular variable is used.⁸ This table plots the raw importance scores for the parameter for weather (x) and nonweather (N) event classes, the Mean Decrease Gini and Mean Decrease Accuracy, which in addition to the OOB error estimate provides us with an estimate of the performance of a variable.

Table 2 Example variable importance measures

Variable	N	x	Mean decrease accuracy	Mean decrease Gini
light_cond	11.840187	23.8492606	23.54195	78.80947
wx_temp	56.976605	-8.3942168	55.818573	377.49762
wx_dew_pt	48.086505	-15.8939562	41.730455	375.30625
wind_dir_deg	1.173042	2.093035	2.136593	378.12715
wind_vel_kts	29.285381	41.6917538	45.632717	338.2876
altimeter	18.951875	-0.4434568	16.739668	443.92041

The raw variable importance scores provide some scores showing how much each variable contributed to the classification of that variable. The Mean Decrease Gini measure is the “sum of all decreases in the forest due to a given variable, normalized by the number of trees”⁹ and gives a total decrease in node impurities from splitting on that variable. When interpreting the Mean Decrease Gini, the larger the value, the purer the variable. The Mean Decrease Accuracy value will use the OOB

samples that were omitted in bootstrapping to determine the importance of the variable as it contributes to the OOB error. The higher the Mean Decrease Accuracy, the more important the variable.⁷ In Table 2, we can see each variable's performance for classification of each nonweather and weather-related classes, the mean decrease accuracy and mean decrease Gini for each variable. In Table 2, we can see that the Mean Decrease Accuracy is the highest for wx_temp, indicating that it performed well with few contributions to the error rate. The Mean Decrease Gini is the highest for altimeter, indicating that it is the purest predictor.

3. Results and Discussion

As described in Section 2.3, the results depend primarily on the OOB error rates, Mean Decrease Gini, and Mean Decrease Accuracy. For each of the data sets, the data have been provided in the Appendix. Here, we broadly discuss some of the final output used to interpret the results. As discussed previously, the higher the Mean Decrease Gini and Mean Decrease Accuracy, the more accurate the results. The OOB error rates allow for a mean prediction error on the bootstrap training samples.

In this report, we compared the following parameters:

- Light conditions (light_cond)
- Temperature (wx_temp)
- Dew point (wx_dew_pt)
- Wind direction (wind_dir_deg)
- Wind velocity (wind_vel_kts)
- Wind gusts (gust_kts)
- Altimeter

The results were divided by the adjustments made in the parameters, as this would significantly change the number of events available for classification. As shown in Table 3, when gusts were included, there were only 1138 airplane events that reported gusts. Meanwhile, when gusts were excluded as a required parameter we had more-complete data records to choose from, resulting in 5590 events to consider. As demonstrated in Table 3, the OOB error rates dropped significantly when more data became available, although they still had approximately 24% error rates.

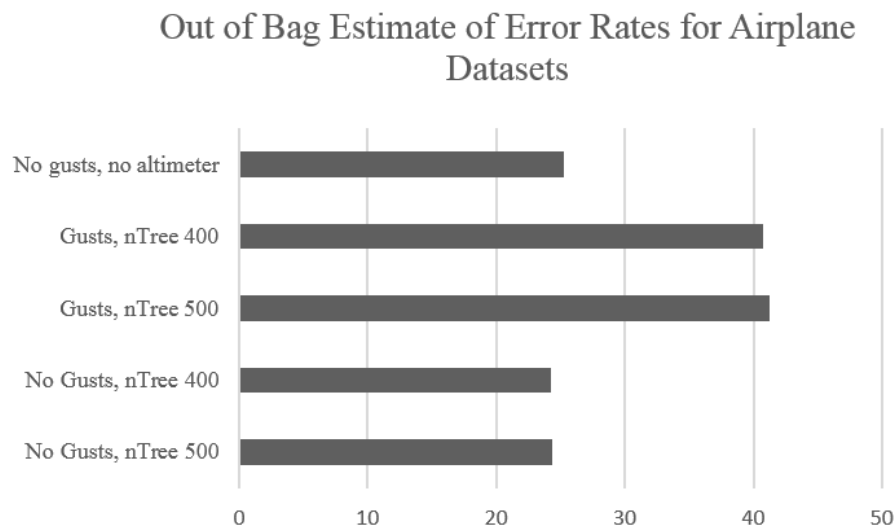
Table 3 Summary of results on airplane data set

Parameter	No. of events	No. of trees	OOB error rates (%)	MVPs
No gusts	5590	400	24.23	Temperature, altimeter
No gusts	5590	500	24.39	Temperature, altimeter
Include gusts	1138	400	40.72	Gusts, altimeter
Include gusts	1138	500	41.42	Gusts, altimeter
No gusts, no altimeter	5590	500	25.28	Temperature, wind direction

Notes: MVP = most valuable predictor; OOB = out of bag.

In the results shown in Table 3, the MVPs often resulted in temperature and altimeter when gusts were excluded. When gusts were included, altimeter and gusts are the MVPs. However, as previously noted, the OOB error rate was also much higher. For all of these data sets, attempts to use both 400 and 500 trees were tried, and this resulted in a small decrease in OOB error rates.

In all of these outcomes, we saw that the outstanding MVP was the altimeter. While this is an important factor, we wanted to see, with altimeter removed, if another environmental parameter would result and what impact this would have on the error rates. For this case, we chose to stay with 500 trees and the largest data set without gusts. Once altimeter was removed, the OOB error rates increased slightly, but the wind direction replaced altimeter as an MVP. Figure 4 shows the performance of each of the OOB error rates for each of the airplane data sets. As previously discussed, the increase in OOB error rates is especially notable in the data sets including gusts, while the other data sets perform relatively well.

**Fig. 4 OOB estimate of error rates for airplane data sets**

With the helicopter data set there was a significant reduction in the number of available events for classification, as shown in Table 4. Without gusts, a maximum of 628 records was found and only 128 records including gusts. While the error rates did relatively well (around 20%) with 628 records, they rose 40%–42% with 128 records. For the no-gust data sets, the temperature and dew point remained the MVPs. However, when gusts were included, light conditions and wind direction became the MVPs instead, keeping in mind the increase in OOB errors (as shown in Fig. 5) and the striking reduction of records available.

While altimeter was not a factor in these MVPs, we attempted to try the same parameter changes performed on the airplane data sets and observe the MVPs and error rates when altimeter was no longer a factor, which resulted in dew point being a clear MVP and a negligible difference in the OOB error rates (Fig. 5). Meanwhile, we did observe light conditions to be a factor when gusts were included; so when excluded, temperature and dew point became MVPs while also seeing a negligible difference in the error rates.

Table 4 Summary of results on helicopter data set

Parameter	No. of trees	No. of events	OOB error rates (%)	MVPs
No gusts	400	628	20.57	Temperature, dew point
No gusts	500	628	20.1	Temperature, dew point
Include gusts	400	128	42.52	Light conditions, wind direction
Include gusts	500	128	40.94	Light conditions, wind direction
No gusts, no altimeter	500	628	20.41	Dew point
No gusts, no altimeter, no light condition	500	628	19.94	Temperature, dew point

Out of Bag Estimate of Error Rates for Helicopter Datasets

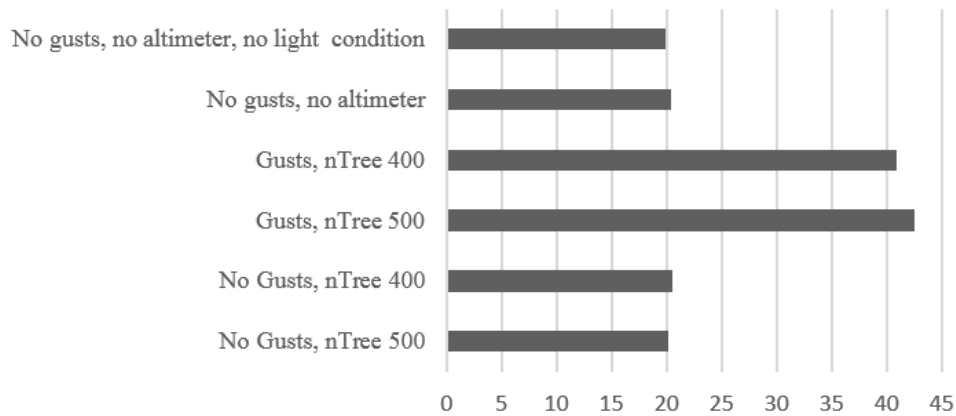


Fig. 5 OOB estimate of error rates for helicopter data sets

4. Conclusions

Across the most accurate data sets, we consistently saw temperature and dew point as factors that could help classify these data sets. However, in all cases, the error rate was significant, typically around 20% in the best cases. While this does demonstrate some interesting results that could contribute to future models, it does not necessarily indicate that this is an accurate outcome and that a model should be derived from the results. When the considerations of safety are a concern, a low error rate is a top priority.

What we can derive from this work is the possibility that should more data become available and more-consistent weather parameters be reported with aircraft incidents, the better the results will be, and a reliable model can be derived from the results. The best scenarios included over 4000 records, while the highest error rates included a fraction of that data.

Future work is required to develop an accurate model for forecasting aircraft incidents. This might include finding larger and consistent data sets that report weather in as many weather parameters as possible. An attempt to supplement the data used in this report could include mining and integration of data from weather stations nearest the incident sites/dates/times in the NTSB records including parameters that were otherwise missing. Additionally, using random forest's regression methods might be useful in an attempt to forecast the events rather than classify them.

5. References

1. National Oceanic and Atmospheric Administration (NOAA). Billion-dollar weather and climate disasters. Asheville (NC): National Centers for Environmental Information (NCEI), NOAA; 2015 May 27 [accessed 2018 Feb 22]. <http://www.ncdc.noaa.gov/billions>.
2. Crawford B. Own the weather: flying in degraded visual environments. Washington (DC): Department of the Army (US); 2015 Jan 19 [accessed 2018 Feb 22]. https://www.army.mil/article/140648/Own_the_Weather_Flying_in_Degraded_Visual_Environments.
3. Judson J. After blimp broke free and crashed: JLENS program hangs by a thread; 2015 Oct 30 [accessed 2017 Apr 19]. <http://defensenews.com/home/2015/10/30/after-blimp-broke-free-and-crashed-jlens-program-hangs-by-a-thread/>.
4. National Transportation Safety Board (NTSB). Aviation data dictionary; 2012 Jan 21 [accessed 2017 May 1]. https://www.nts.gov/_layouts/ntsb.aviation/AviationDownloadDataDictionary.aspx.
5. Liaw A. Package “randomForest”: the comprehensive R archive network; 2015 Oct [accessed 2017 Sep 19]. <https://cran.r-project.org/web/packages/randomForest/randomForest.pdf>.
6. Johnson JO, White Sands Missile Range, NM. Human vulnerability forecasting. Personal communication; 2016 May.
7. Han H, Guo X, Yu H. Variable selection using mean decrease accuracy and mean decrease Gini based on random forest. Piscataway (NJ): The Institute of Electrical and Electronics Engineers (IEEE); 2016. Paper No.: 978-1-4673-9904-3.
8. Liaw A, Wiener M. Classification and regression by random forest. R News. 2002;2/3:18–22.
9. Breiman L, Cutler A. Manual: setting up, using, and understanding random forests, v4.0. Berkeley (CA): University of California Berkeley; 2003.

INTENTIONALLY LEFT BLANK.

Appendix. Database Queries and Random Forest Output

Helicopter Access Query:

```
SELECT DISTINCT FINDINGS.EV_ID, FINDINGS.FINDING_DESCRIPTION,  
EVENTS.LIGHT_COND, EVENTS.WX_TEMP, EVENTS.WX_DEW_PT,  
EVENTS.WIND_DIR_DEG, EVENTS.WIND_VEL_KTS, EVENTS.GUST_KTS,  
EVENTS.ALTIMETER  
FROM FINDINGS, EVENTS, AIRCRAFT  
WHERE (((FINDINGS.EV_ID)=[AIRCRAFT].[EV_ID] AND  
(FINDINGS.EV_ID)=[EVENTS].[EV_ID]) AND ((AIRCRAFT.ACFT_CATEGORY) LIKE  
'HELI'));
```

Airplane Access Query:

```
SELECT DISTINCT FINDINGS.EV_ID, FINDINGS.FINDING_DESCRIPTION,  
EVENTS.LIGHT_COND, EVENTS.WX_TEMP, EVENTS.WX_DEW_PT,  
EVENTS.WIND_DIR_DEG, EVENTS.WIND_VEL_KTS, EVENTS.GUST_KTS,  
EVENTS.ALTIMETER  
FROM FINDINGS, EVENTS, AIRCRAFT  
WHERE (((FINDINGS.EV_ID)=[AIRCRAFT].[EV_ID] AND  
(FINDINGS.EV_ID)=[EVENTS].[EV_ID]) AND ((AIRCRAFT.ACFT_CATEGORY) LIKE  
'AIR*'));
```

Detailed Results from Random Forest Runs

Airplanes

No Gusts – ntree 500

Random Forest output:

Type of random forest: classification

Number of trees: 500

No. of variables tried at each split: 2

Out-of-bag (OOB) estimate of error rate: 24.39%

Table A-1 Confusion matrix airplanes, ntree 500, no gusts

	N	x	class.error
N	3963	260	0.06156761
x	1103	263	0.80746706

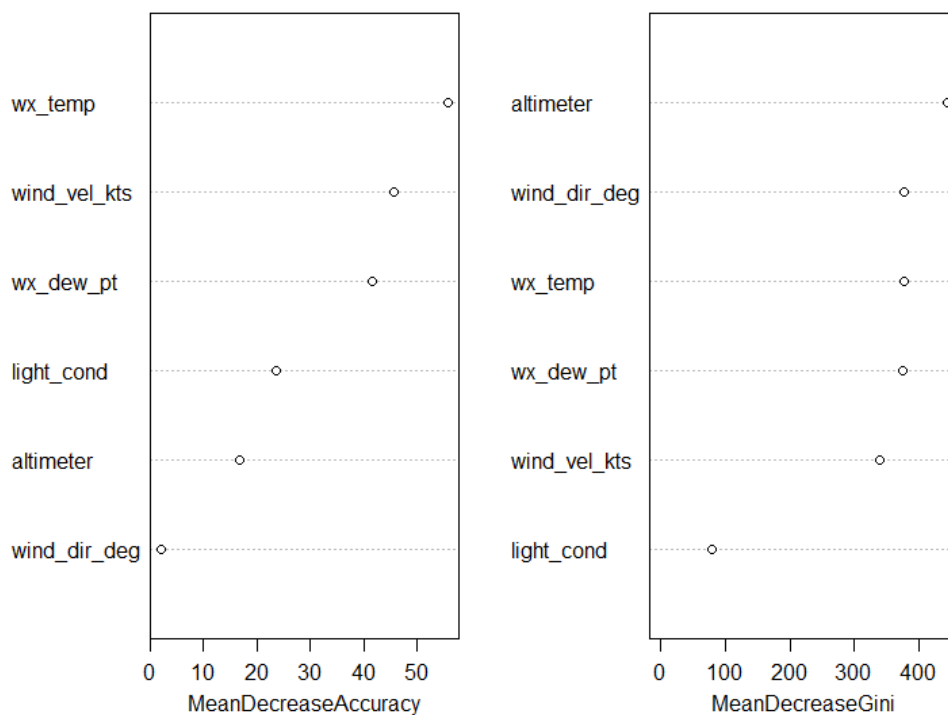


Fig. A-1 Variable importance plot: airplanes, ntree 500, no gusts

Table A-2 Variable importance measures: airplanes, ntree 500, no gusts

Variable	N	x	Mean decrease accuracy	Mean decrease Gini
light_cond	11.840187	23.8492606	23.54195	78.80947
wx_temp	56.976605	-8.3942168	55.818573	377.49762
wx_dew_pt	48.086505	-15.8939562	41.730455	375.30625
wind_dir_deg	1.173042	2.093035	2.136593	378.12715
wind_vel_kts	29.285381	41.6917538	45.632717	338.2876
altimeter	18.951875	-0.4434568	16.739668	443.92041

No Gusts – ntree 400

Random Forest output:
Type of random forest: classification
Number of trees: 400
No. of variables tried at each split: 2
OOB estimate of error rate: 24.23%

Table A-3 Confusion matrix for airplanes, ntree 400, no gusts

	<i>N</i>	<i>x</i>	class.error
N	3970	253	0.05991002
x	1101	265	0.80600293

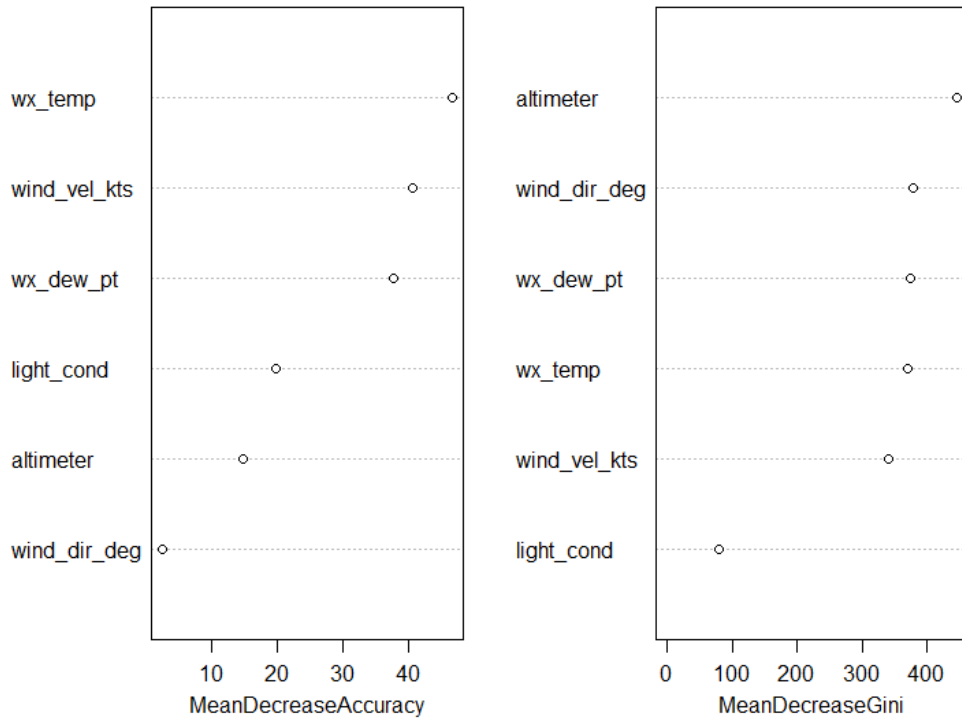


Fig. A-2 Variable importance plot: airplanes, ntree 400, no gusts

Table A-4 Variable importance measures: airplanes, ntree 400, no gusts

Variable	N	x	Mean decrease accuracy	Mean decrease Gini
light_cond	10.0116	21.8751953	19.742184	80.18966
wx_temp	47.5407	-7.8439718	46.748625	371.4269
wx_dew_pt	42.20958	-13.4995938	37.799126	373.95258
wind_dir_deg	2.66869	0.3804706	2.442168	380.0827
wind_vel_kts	26.74702	32.0925775	40.623684	341.41796
altimeter	17.8817	-3.1895253	14.674056	446.21703

Gusts – ntree 500

Random Forest output:

Type of random forest: classification

Number of trees: 500

No. of variables tried at each split: 2

OOB estimate of error rate: 41.42%

Table A-5 Confusion matrix airplanes, ntree 500, gusts

	<i>N</i>	<i>x</i>	class.error
N	452	182	0.2870662
x	289	214	0.5745527

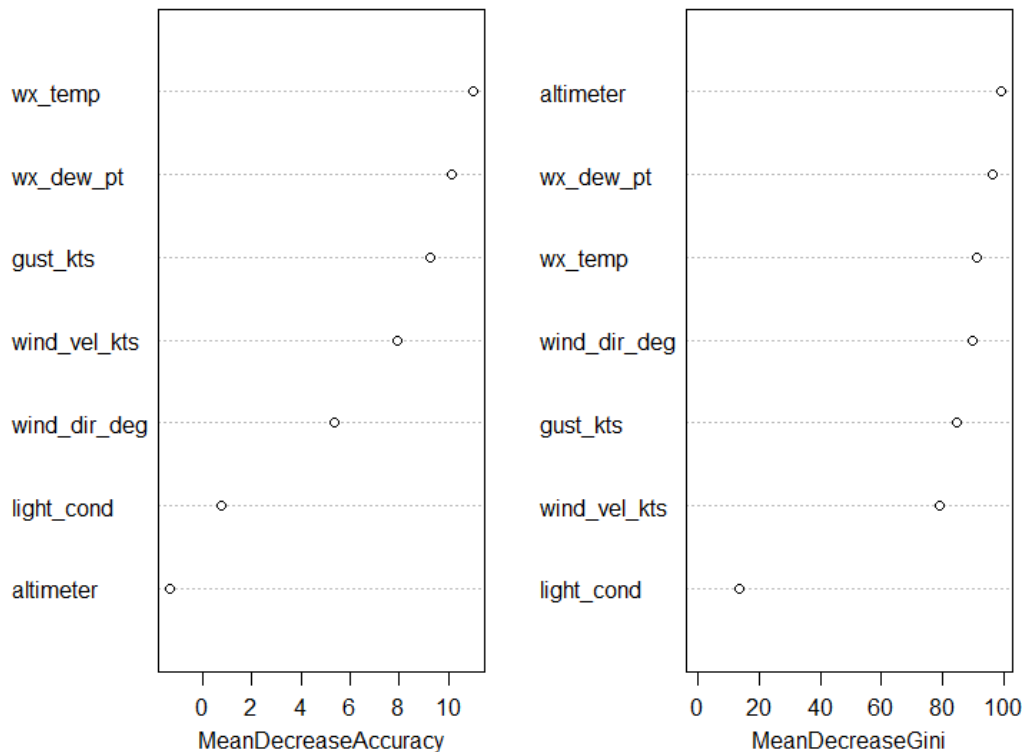


Fig. A-3 Variable importance plot: airplanes, ntree 500, gusts

Table A-6 Variable importance measures: airplanes, ntree 500, gusts

Variable	N	x	Mean decrease accuracy	Mean decrease Gini
light_cond	2.233281	-1.2191109	0.7422473	13.4635
wx_temp	11.50686	1.4360709	10.9949025	91.33929
wx_dew_pt	13.306476	-0.3513214	10.1174796	96.1576
wind_dir_deg	8.347022	-1.8519558	5.3281705	89.57317
wind_vel_kts	7.899476	1.8112498	7.8923629	79.14064
altimeter	4.984241	-8.1366573	-1.3194029	99.07035
gust_kts	8.767074	2.7485622	9.2469088	84.6994

Gusts – ntree 400

Random Forest output:

*Type of random forest: classification**Number of trees: 400**No. of variables tried at each split: 2**OOB estimate of error rate: 40.72%***Table A-7 Confusion matrix, airplanes, ntree 400, gusts**

	N	x	class.error
N	453	181	0.285489
x	282	221	0.5606362

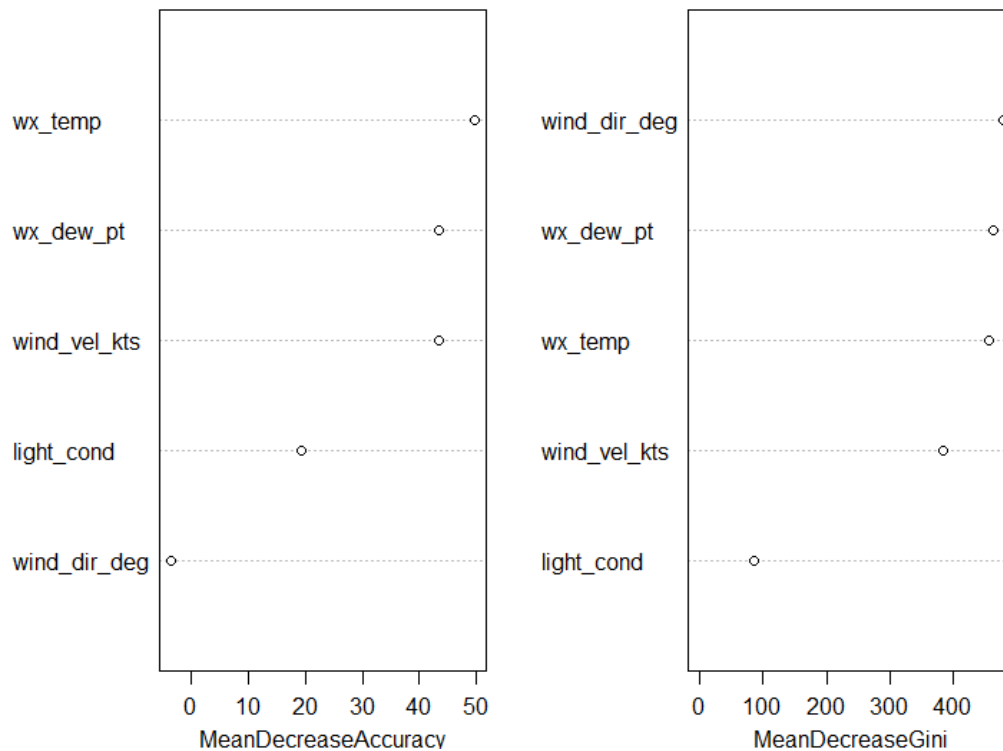


Fig. A-4 Variable importance plot:-airplanes, ntree 400, gusts

Table A-8 Variable importance measures: airplanes, ntree 400, gusts

Variable	N	x	Mean decrease accuracy	Mean decrease Gini
light_cond	3.572651	-0.1039458	2.312476	13.82871
wx_temp	11.198215	-1.092739	9.084722	92.49765
wx_dew_pt	11.084761	-2.2840685	8.22982	94.22073
wind_dir_deg	5.742023	-2.505355	2.552762	89.67124
wind_vel_kts	6.75407	0.5237188	6.005231	77.31019
altimeter	3.999004	-7.9153185	-1.808045	101.00583
gust_kts	8.819424	2.3902855	9.438109	84.37634

No Gusts - No Altimeter

Random Forest output:

Type of random forest: classification

Number of trees: 500

No. of variables tried at each split: 2

OOB estimate of error rate: 25.28%

Table A-9 Confusion matrix, airplanes, ntree 500, no gusts, no altimeter

	<i>N</i>	<i>x</i>	class.error
N	3873	350	0.08287947
x	1063	303	0.77818448

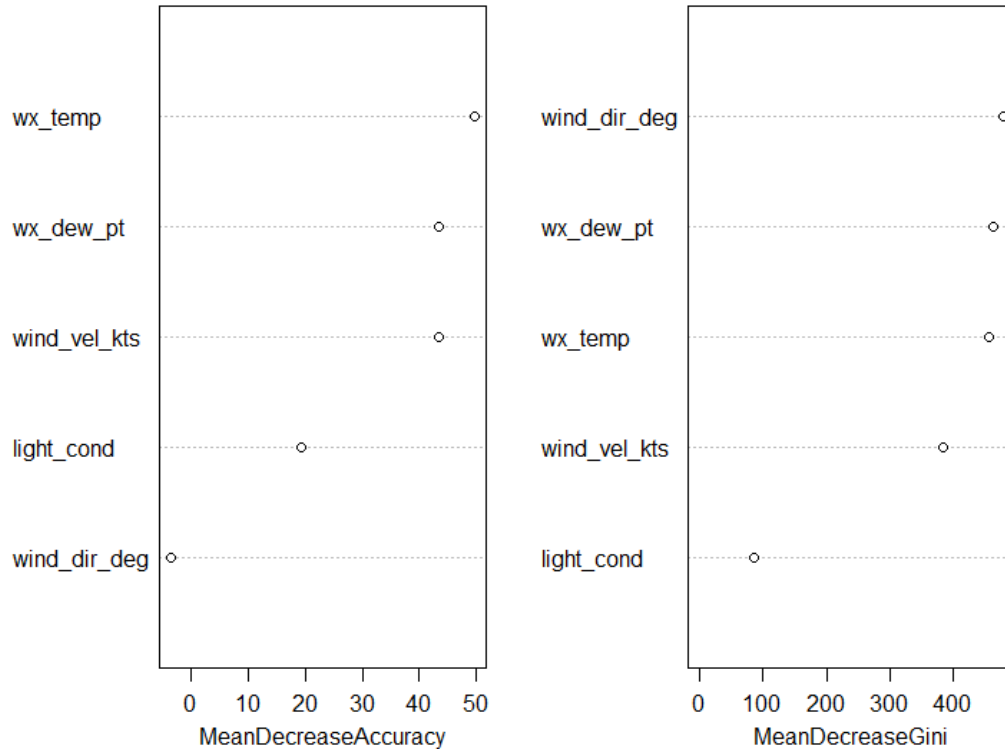


Fig. A-5 Variable importance plot: airplanes, ntree 500, gusts, no altimeter

Table A-10 Variable importance measures: airplanes, ntree 500, no gusts, no altimeter

Variable	N	x	Mean decrease accuracy	Mean decrease Gini
light_cond	10.75	18.4658838	19.401291	85.39137
wx_temp	49.38828	-3.1704462	49.614706	456.58398
wx_dew_pt	46.49761	-13.7953901	43.423776	463.79049
wind_dir_deg	-4.27994	0.4630378	-3.376186	479.53358
wind_vel_kts	25.11055	41.233236	43.407036	384.0949

Helicopters

No Gusts – ntree 500

Random Forest output:

Type of random forest: classification

Number of trees: 500

No. of variables tried at each split: 2

OOB estimate of error rate: 20.1%

Table A-11 Confusion matrix, helicopters, ntree 500, no gusts

	<i>N</i>	<i>x</i>	class.error
N	489	20	0.03929273
x	106	12	0.89830508

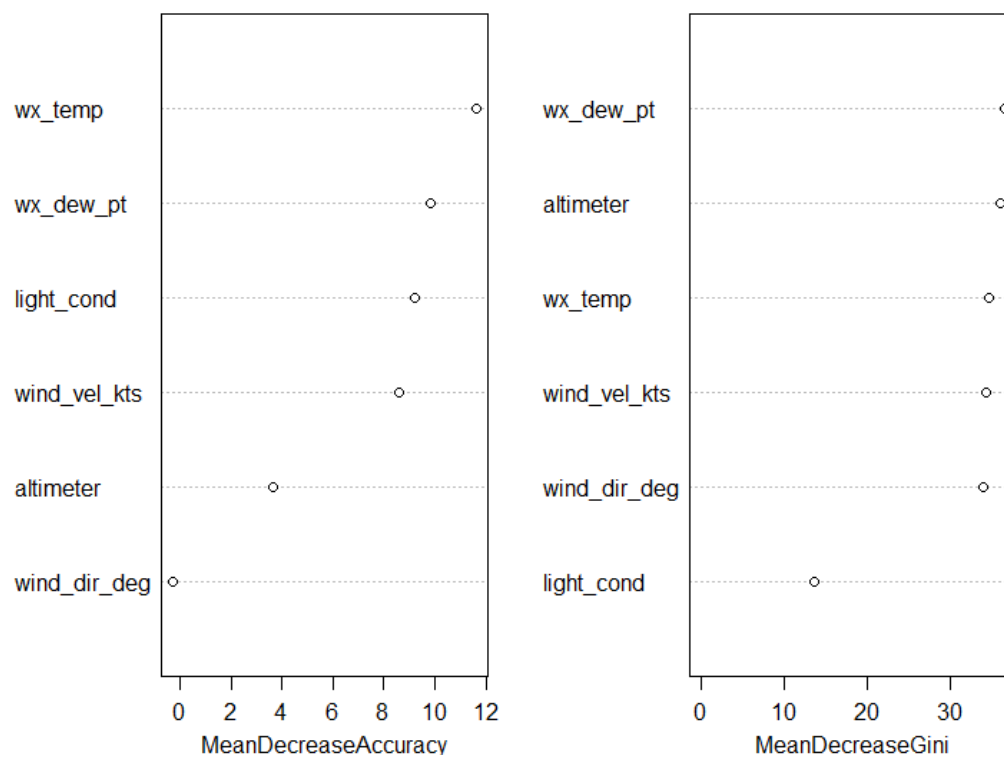


Fig. A-6 Variable importance measures: helicopter, ntree 500, no gusts

Table A-12 Variable importance plot: helicopter, ntree 500, gusts

Variable	N	x	Mean decrease accuracy	Mean decrease Gini
light_cond	1.91156587	12.975265	8.4309547	13.02388
wx_temp	13.32815904	-1.692774	12.3363263	34.91327
wx_dew_pt	15.00538148	-8.915313	10.8841943	36.77802
wind_dir_deg	0.02346716	-1.89521	-0.8780432	33.71206
wind_vel_kts	8.2319426	7.010469	10.5067365	33.96037
altimeter	6.51944018	-5.088066	3.8253575	35.93895

No Gusts – ntree 400

Type of random forest: classification

Number of trees: 400

No. of variables tried at each split: 2

OOB estimate of error rate: 20.57

Table A-13 Confusion matrix for helicopter, no gusts, ntree 400

	N	x	class.error
N	487	22	0.043222
x	107	11	0.9067797

Table A-14 Variable importance measures: helicopter, ntree 400, no gusts

Variable	N	x	Mean decrease accuracy	Mean decrease Gini
light_cond	3.761427	11.0426758	9.1891123	13.5171
wx_temp	12.4511174	-2.2359166	11.6197874	34.62552
wx_dew_pt	13.1803534	-7.9901741	9.8062402	36.53487
wind_dir_deg	-0.1023798	-0.6490411	-0.2731138	34.02029
wind_vel_kts	7.3861387	4.6708483	8.5827052	34.25989
altimeter	6.1154741	-4.1367783	3.627002	35.96782

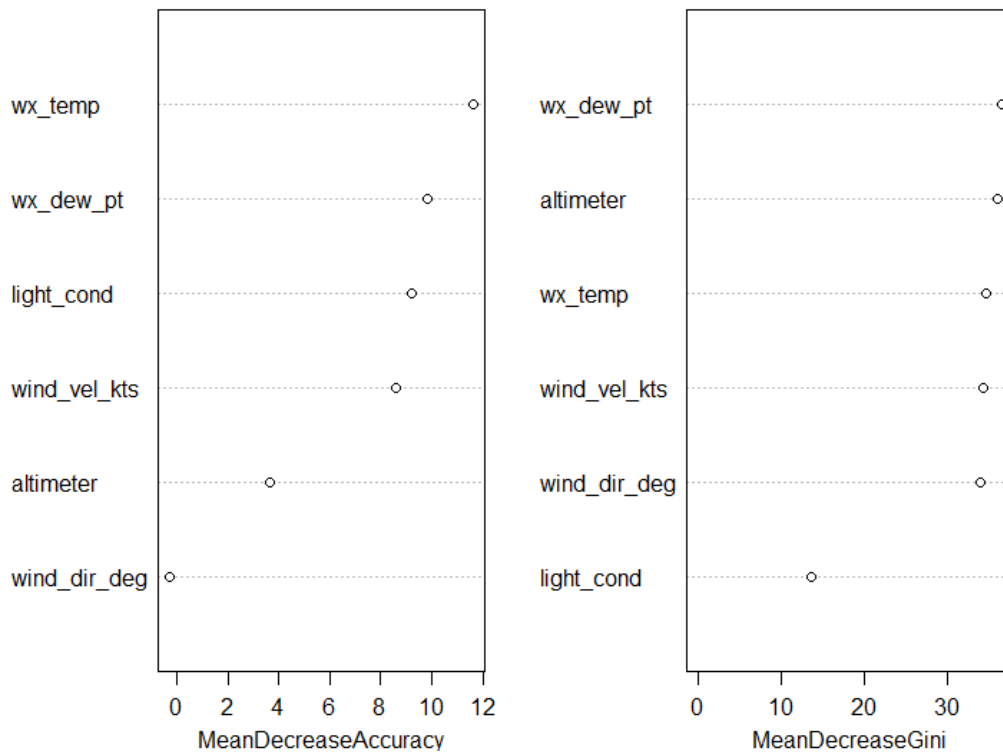


Fig. A-7 Variable importance plot: helicopter, ntree 400, no gusts

Gusts – ntree 500

Random Forest Output:

Type of random forest: classification

Number of trees: 500

No. of variables tried at each split: 2

OOB estimate of error rate: 40.94%

Table A-15 Confusion matrix for helicopter, gusts, ntree 500

	<i>N</i>	<i>x</i>	class.error
N	67	14	0.1728395
x	38	8	0.826087

Table A-16 Importance measures for helicopter, gusts, ntree 500

Variable	N	x	Mean decrease accuracy	Mean decrease Gini
light_cond	4.5045482	5.2752305	6.39485335	3.104077
wx_temp	0.6844501	-0.933441	0.03292168	9.075437
wx_dew_pt	1.501178	-4.8069218	-1.68761907	8.300828
wind_dir_deg	-1.4421245	0.3443367	-0.92445561	10.70617
wind_vel_kts	0.2616115	0.343606	0.29784493	8.454298
altimeter	-0.5181441	-4.1327118	-2.94821493	9.633624
gust_kts	0.4900902	2.8823347	1.95880135	8.715406

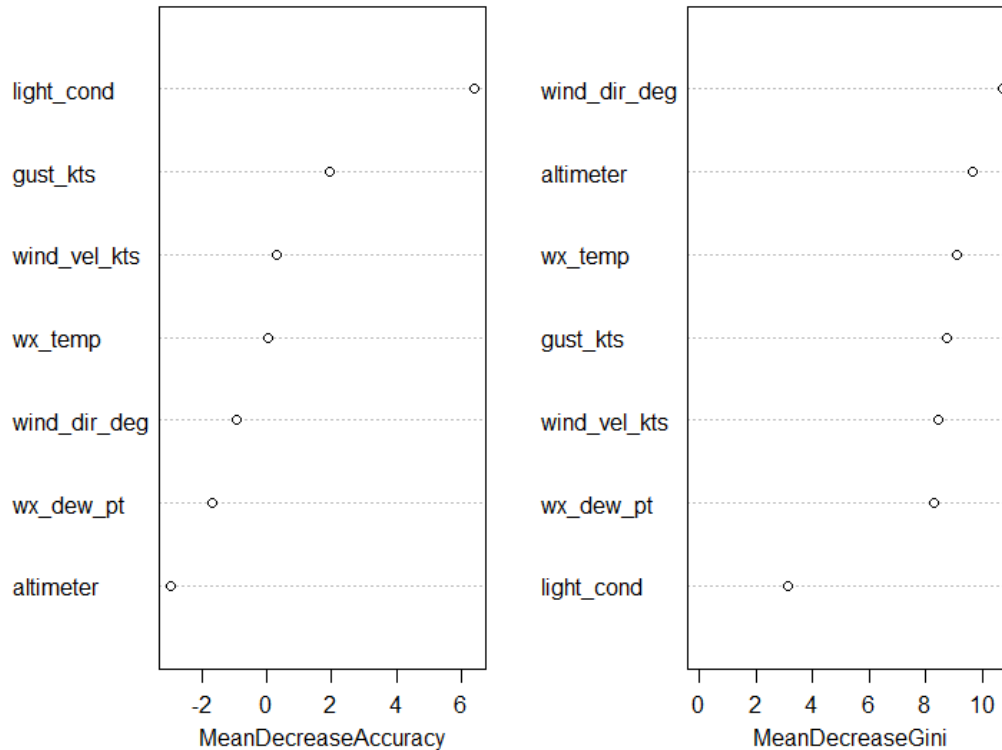


Fig. A-8 Variable importance plot: helicopter, gusts, ntree 500

Gusts – ntree 400

Type of random forest: classification

Number of trees: 400

No. of variables tried at each split: 2

OOB estimate of error rate: 42.52%

Table A-17 Confusion matrix for helicopter, gusts, ntree 400

	<i>N</i>	<i>x</i>	class.error
N	62	19	0.2345679
x	35	11	0.7608696

Table A-18 Importance measures for helicopter, gusts, ntree 400

Variable	N	x	Mean decrease accuracy	Mean decrease Gini
light_cond	1.30499145	7.58419	5.6425638	3.040114
wx_temp	-0.07109502	0.8202233	0.4865997	9.281376
wx_dew_pt	1.17175679	-5.2582878	-1.9722996	8.58466
wind_dir_deg	-1.95085433	1.8220211	-0.3639487	10.759474
wind_vel_kts	-0.3350203	-0.7978229	-0.7744061	8.203409
altimeter	-1.96459161	-3.9434259	-4.1299063	9.362583
gust_kts	0.06931692	2.4782635	1.2987592	8.66372

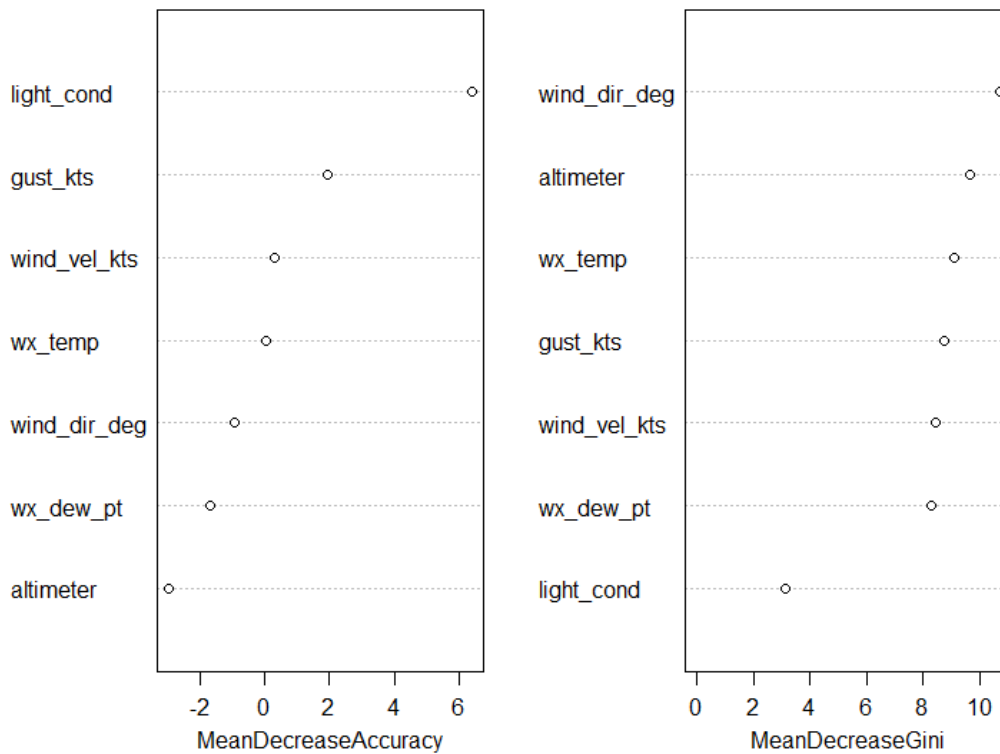


Fig. A-9 Variable importance plot: helicopter, gusts, ntree 400

No Altimeter

Number of trees: 500

No. of variables tried at each split: 2

OOB estimate of error rate: 20.41%

Table A-19 Confusion matrix for helicopter, no gusts, no altimeter, ntree 500

	<i>N</i>	<i>x</i>	class.error
N	485	24	0.04715128
x	104	14	0.88135593

Table A-20 Variable importance measures for helicopter, no gusts, no altimeter, ntree 500

Variable	N	x	Mean decrease accuracy	Mean decrease Gini
light_cond	2.1342769	13.2661324	8.5855959	13.80308
wx_temp	13.0161326	0.6226949	13.2557421	42.05545
wx_dew_pt	19.1026043	-9.0211285	14.7868566	46.55484
wind_dir_deg	-0.2575318	-1.2728565	-0.8292119	43.8914
wind_vel_kts	11.1414579	9.3122348	14.2896222	40.22386

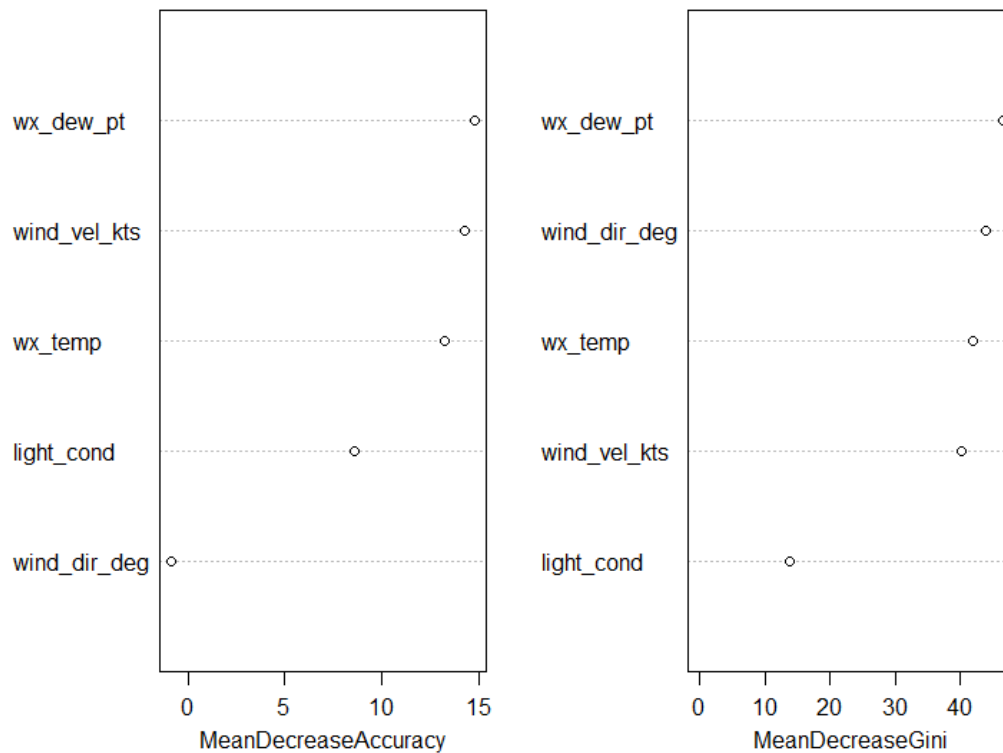


Fig. A-10 Variable importance plot for helicopter, no gusts, no altimeter, ntree 500

List of Symbols, Abbreviations, and Acronyms

ID	identification
JLENS	Joint Land Attack Cruise Missile Defense Elevated Netted Sensor System
MVP	most valuable predictor
NTSB	National Transportation Safety Board
OOB	out of bag

1 DEFENSE TECHNICAL
(PDF) INFORMATION CTR
DTIC OCA

2 DIR ARL
(PDF) IMAL HRA
RECORDS MGMT
RDRL DCL
TECH LIB

1 GOVT PRINTG OFC
(PDF) A MALHOTRA

1 DIR ARL
(PDF) RDRL CIE D
Y RABY